

# PROYECTO BIG DATA Y VISUALIZACIÓN

Juan Esteban Echeverri Peña, Jhonny Steven Ibarbo Herrera, Brayan Urdinola Sanchez

*Ingeniería Electrónica, Universidad Autónoma de Manizales - UAM*  
Manizales, Colombia

24 de Mayo de 2023

---

**Resumen:** El proyecto consta de dos partes fundamentales. En la primera parte, se emplea la inteligencia artificial de Whisper para convertir audio en texto. Esto se logra cargando el modelo de Whisper y realizando la transcripción del audio utilizando dicho modelo. El resultado se almacena en un DataFrame de Pandas y se convierte posteriormente en un DataFrame de Spark. A partir de este DataFrame de Spark, se realizan consultas SQL para obtener un resultado específico. Finalmente, el resultado se guarda en una variable de tipo cadena para su posterior uso.

En la segunda parte del proyecto, se lleva a cabo la creación de bases de datos para tres géneros musicales diferentes: Rap, Hip Hop y Electrónica. Se ingresan tres canciones en cada base de datos y se unifican en una única variable. A continuación, se realiza un proceso de filtrado para eliminar palabras y caracteres no deseados en la variable que contiene las canciones. Luego, se utiliza Spark para contar las palabras más y menos utilizadas en cada género musical, así como en la canción que fue transcrita en la primera parte del proyecto. Por último, se generan gráficos que representan las palabras más y menos utilizadas en cada género musical.

El proceso general del proyecto se organiza en un menú interactivo que permite al usuario seleccionar las opciones deseadas. El código incluye la función "inicio()" que implementa este menú. En él, el usuario puede elegir trabajar con texto o audio, y realizar diferentes acciones como ingresar la letra de una canción, visualizar análisis existentes, o salir del programa. Además, se proporciona la opción de analizar y graficar las palabras más y menos utilizadas en los géneros musicales de la base de datos, así como en la canción transcrita.

**Palabras clave:** Transcriptor de audio , Inteligencia artificial, DataFrame, Análisis de palabras, Spark, MapReduce.

---

## I. OBJETIVO

- Desarrollar un sistema capaz de transcribir audios a texto y utilizar técnicas de procesamiento distribuido, como MapReduce o el framework de Apache Spark, para extraer las palabras más y menos utilizadas tanto de bases de datos predefinidas como del audio o texto ingresado.

## II. INTRODUCCIÓN

En el campo de los audiovisuales, donde la cantidad de datos disponibles es cada vez mayor, se ha vuelto evidente la necesidad de aprovechar las técnicas y herramientas de Big Data para obtener un mayor entendimiento del entorno en el que nos desenvolvemos. La disponibilidad de datos en diversas formas, como textos, audios y bases de datos, plantea

el desafío de procesar y analizar esta vasta cantidad de información de manera eficiente y efectiva.

En este contexto, el presente proyecto se centra en la transcripción de audios y el análisis de palabras de archivos de audio. La transcripción de audios permite convertir información hablada en texto, lo cual facilita su posterior procesamiento y análisis. Por otro lado, el análisis de palabras nos brinda información valiosa sobre los patrones lingüísticos y el contenido temático presente en los datos.

Para abordar estos desafíos, se emplean técnicas de Big Data, en particular MapReduce y el framework de Apache Spark. Estas herramientas son ampliamente reconocidas por su capacidad de procesar grandes volúmenes de datos de manera distribuida y escalable. MapReduce divide tareas complejas en etapas más pequeñas y las distribuye en clústeres de computadoras, lo que acelera el procesamiento y permite manejar conjuntos de datos masivos. Por su parte, Apache Spark es un framework de procesamiento distribuido que facilita el análisis de datos en tiempo real y ofrece una amplia gama de herramientas y bibliotecas para el procesamiento y análisis de datos a gran escala.

Una vez transcritos los audios, el texto resultante se almacena en DataFrames, estructuras de datos tabulares que permiten organizar y manipular los datos de manera eficiente. En este proyecto, se utilizan DataFrames de Pandas y Apache Spark para almacenar los resultados de la transcripción y el análisis de palabras. Estas tablas proporcionan una representación estructurada de los datos, lo que facilita su visualización y consulta posterior.

Además, el uso de MapReduce y Apache Spark ayuda significativamente en la organización y extracción de información que

se muestra en las gráficas de barras. Estas herramientas permiten realizar cálculos distribuidos y paralelos sobre los conjuntos de datos, lo que acelera el procesamiento y análisis de las palabras extraídas. De esta manera, se obtienen resultados más rápidos y precisos en las gráficas de barras, lo que facilita la identificación de las palabras más y menos utilizadas en las diferentes categorías analizadas.

### III. METODOLOGÍA

Inicialmente se destaca que los 3 integrantes, Juan Esteban Echeverri Peña, Jhonny Steven Ibabo Herrera y Brayan Urdinola Sanchez; participaron en cada una de las etapas empleadas a continuación:

#### 1. Definición del objetivo del proyecto:

Consiste en desarrollar un sistema capaz de transcribir audios a texto y analizar las palabras más y menos utilizadas en el campo seleccionado.

#### 2. Investigación y selección de herramientas y técnicas:

Se llevó a cabo una investigación exhaustiva sobre las herramientas y técnicas de Big Data disponibles para la transcripción de audio y el análisis de palabras. Se evaluaron diferentes opciones y se seleccionaron MapReduce y Apache Spark como los frameworks principales a utilizar.

#### 3. Diseño del flujo de trabajo:

Se establecieron las etapas principales, como la carga del modelo de transcripción de audio, la conversión de datos a DataFrames de Pandas y Apache Spark, las consultas SQL, el filtrado de palabras no deseadas y la generación de gráficas de barras.

#### 4. Implementación del código:



Posteriormente se realizó una consulta mediante SQL, en este caso se solicitó los datos 'start', 'end' y 'text':

start	end	text
0.0	122.55	Tengo un número, Lore No se pa' qué, si me lo sé en memoria
122.55	136.34000000000001	Me hiciste daño y así te sentí
136.34	134.8	Y aunque sé que un día te voy a olvidar Aún no lo hago, es complicado
134.8	141.28	Lo traidito me gusta recordar Ando manejando por las calles que me besaste
141.28	149.36	O viendo las canciones que un día me dedicaste Te dije que volverías pero eso no se pide
149.36	155.64	Mejor le pido a Dios que me cuide Porque ando manejando por las calles que me besaste
155.64	164.4	O viendo las canciones que un día me dedicaste Te dije que volverías pero eso no se pide
164.4	171.6	Mejor le pido a Dios que me cuide Que me pida otra que se parezca a ti
171.6	179.68	No quiero estar como hice por ti Ojalá te enamore
179.68	188.72000000000001	Lo mismo que me hiciste a mí Tu me enseñaste a no near a cualquiera
188.72000000000001	195.56	Y cambié como yo quiero Me daban tres en una relación de dos
195.56	199.6	No te perdono, pídele perdón a Dios Dije que te olvidé y la verdad es que yo
199.6	1113.52	Ando manejando por las calles que me besaste O viendo las canciones que un día me dedicaste
1113.52	1121.1999999999999	Te dije que volverías pero eso no se pide Mejor le pido a Dios que me cuide
1121.1999999999999	1129.6	Te dije que volverías pero eso no se pide Mejor le pido a Dios que me cuide

Figura 4: Consulta por SQL

Finalmente, se realiza la transferencia de los datos del DataFrame a una variable llamada "audiocan". Para lograr esto, se utiliza la función collect(), que permite recopilar los resultados del DataFrame. A su vez, se utiliza la expresión "row.text" para acceder al valor de la columna "text" de cada fila del DataFrame.

Es importante destacar que la variable "audiocan" se almacena como una lista. Sin embargo, para poder aplicar el filtrado y el graficado requeridos en este proyecto, es necesario convertir esta lista en una cadena de texto. Por esta razón, se optó por modificar la lista "audiocan" y transformarla en una cadena llamada "audiocan\_str" utilizando el método "join(audiocan)" junto con un espacio en blanco como separador. Este proceso permite obtener una representación en formato de cadena de los datos del audio o canción transcrita, facilitando así su posterior manipulación y análisis.

```
1 audiocan_str = ' '.join(audiocan) # Pasamos la lista de 'audiocan' a
2 print("audiocan como cadena de texto:", audiocan_str)
```

audiocan como cadena de texto: Tengo un número, Lore No se pa' qué, si me

Figura 5: Conversión de lista a String

## CONJUNTOS DE DATOS:

Se crean conjuntos de datos para tres géneros musicales distintos: Rap, Hip Hop y Electrónica. Adicionalmente se tiene la canción transcrita y una canción adicional que se puede ingresar.

Posteriormente es necesario realizar un filtrado

de palabras y caracteres no deseados en las canciones, esto implica eliminar palabras vacías, signos de puntuación y caracteres especiales que no aportan información relevante para el análisis de las palabras utilizadas en las canciones.

Una vez realizada la limpieza de las canciones, se emplea MapReduce o el framework de Apache Spark para llevar a cabo el conteo de palabras más y menos usadas en cada género musical y en la canción ingresada previamente.

Finalmente se procede a graficar las palabras más y menos usadas en cada género musical. Estas gráficas de barras permiten visualizar de manera clara y concisa cuáles son las palabras más frecuentes y menos frecuentes en cada género, así como en la canción ingresada:

### - Gráfica - Rap:

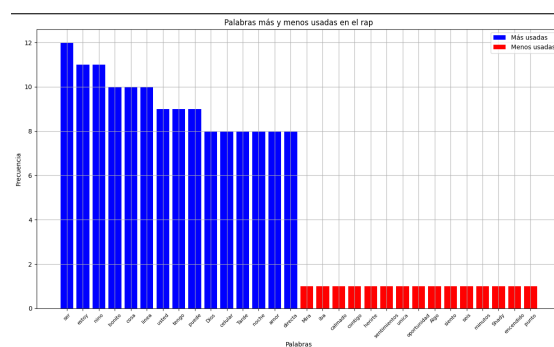


Figura 6: Palabras más y menos usadas en el rap.

Teniendo en cuenta la gráfica anterior, se puede concluir que en el rap, se abordan temas como la realidad social, la lucha, la superación personal, el amor, la desigualdad y la crítica social. Palabras relacionadas con estas temáticas, como "ser, niño, tarde, noche, amor", son utilizadas frecuentemente para transmitir mensajes y narrativas propias del género.



## VI. REFERENCIAS

[1] OpenAI, “whisper” (2023). Disponible en:  
GitHub, recuperado de:  
<https://github.com/openai/whisper>

[2] Anónimo, “Letras de canciones” (2023).  
Disponible en: letras, recuperado de:  
<https://www.letras.com/>